

1 Categorical-categorical interactions

Let's consider two models that correspond to two possible theories about the world. First, $W = \beta_0 + \beta_1 U + \beta_2 F + \epsilon$, where W is wage, U is union status and F is whether one is female, and a second model where $W = \beta_0 + \beta_1 U + \beta_2 F + \beta_3 UF + \epsilon$. **Does the first model here have an interaction?**

Estimate this first model using ordinary least squares; interpret the results. Find the conditional regression equation for men and for women as a function of union status and the conditional regression equation for union and non-union members as a function of gender; predict the values for the four groups the two dummies logically imply exist. To do so, let's use the new CPS 2019 extract I put up (sorry for the many versions; it is a huge file, but it is much better than the GSS for wage information). You can also just use the full file if you've downloaded it from the CEPR website before. Use `wage4`, `female`, and `union`.

Let's see how to verify your hand-calculations¹ using a new command, `margins` (I've shown it briefly before). This must be run right after the regression; it is a post-estimation command. You should also run the regression with Stata's *factor variable syntax* and *interaction syntax* (we'll show the interaction later). `margins` is an extremely flexible, very useful command with an enormous number of options.

```
reg wage4 i.female i.union
margins, at(female=(0 1) union=(0 1))
marginsplot
* Let's transpose the graph
marginsplot, xdimension(union)
```

What do you notice about the graphical presentation?

Now, let's specify our a different model, $Y = \beta_0 + \beta_1 D_1 + \beta_2 D_2 + \beta_3 D_1 D_2 + \epsilon$. **Does this model have an interaction? What will our coefficients tell us?**

Estimate the fully-interacted model. From the model you estimated, calculate “by hand” the regression equation for men and women and union/non-union members, and calculate the predicted values for all four groups.² Check with `margins` and use `marginsplot`; What do you notice about the graphical presentation of the results?

```
reg wage4 i.female##i.union
margins, at(female=(0 1) union=(0 1))
marginsplot
* Note that if you don't want to manually specify
```

¹Your overall equation should be something like $\hat{W} = 27.75 - 4.66F + 4.10U$. Your conditional equation for non-union members is $\hat{W} = 27.75 - 4.66F$; for union members, $\hat{W} = (27.75 + 4.10) - 4.66F$. Your conditional equation for women should be $\hat{W} = (27.75 - 4.66) + 4.10U$; for men, it should be $\hat{W} = 27.75 + 4.10U$. Then, your predicted values can be obtained with plug-and-chug: for non-union men and women, we simply have $\hat{W} = 27.75 - 4.66(0) = 27.75$; $\hat{W} = 27.75 - 4.66(1) = 23.09$. For union men and women, we have $\hat{W} = (27.75 + 4.10) - 4.66(0) = 27.19$; $(27.75 + 4.10) - 4.66(1) = 31.85$.

²The regression equation for each group should be something like this. For men, $\hat{W} = 27.86 + 3.17U$. For women, $\hat{W} = 22.98 + 5.13U$. For non-union members, your equation should be something like $\hat{W} = 27.86 - 4.88F$; for union members, $\hat{W} = 31.03 - 2.93F$. For Your predicted values should be as follows: non-union men, 27.86; union men, 31.03; non-union women, 22.98; union women, 28.1.

```
* the values, just use -levelsof-
levelsof female, local(sexes)
levelsof union, local(unions)
margins, at(female = (`sexes') union = (`unions'))
marginsplot
```

It is useful to make sure that we know what an interaction really mean. It is tempting, for example, to assume that in a case involving race and gender predictors, the thesis of *intersectionality* might suggest the need for an *interaction* effect. One popular rendering of such theories is that “oppression stacks”. **Thought question: in such a phrasing, is an interaction model indicated or not?**³

Whatever your answer to the above question, a clearer example might be useful to keep in mind. One classic and literal example is drug interactions. Some are bad; some are good. However, all share the property that the independent effects of each drug do not simply combine but each amplifies the effect of the other.

1.1 Testing other hypotheses

1.1.1 Re-parameterizing the model

Note that in the above, we did not directly test some potentially-interesting hypotheses: that non-union women get a union benefit and that there is a gender effect among union members. We calculated those regression equations by hand, but they did not show up in our default regression output, and therefore, since we also did not test their significance when calculating by hand, we don’t know if those effects are actually significant.

How can we test those claims? We have the same three options from before: 1) re-parameterizing the model; 2) an F -test; 3) a test of a linear combination of coefficients.

Let’s start with a different parameterization of the model: $W = \beta_0 + \beta_1 U + \beta_2 M + \beta_3 UM + \epsilon$. **Make appropriate dummies and a product term, or just cleverly use Stata’s syntax.** I’ll show the latter below. Recall that the option `nohead` suppresses the “ANOVA table” (sums of squares and tests based on that), which makes it easier to compare slopes. **Thought question: can any part of the sums of squares output change based on what we are doing?**

```
reg wage4 ib1.female##i.union, nohead
```

What does the t -statistic on `union` now correspond to? It tests whether there is a union effect for women. Our baseline model only directly estimated this for men.

Now finally write our model like so... $W = \beta_0 + \beta_1 N + \beta_2 F + \beta_3 NF + \epsilon$. Then estimate.

```
reg wage4 i.female##ib1.union
```

³Of course, there is no one right answer. However, if the person simply means “there is racism against group $A = a$ and sexism against group $B = b$, therefore it is exceptionally difficult to be in group ab ”, this is **not** an interaction model of necessity! Why? The “extra hard” could simply be additive. An interaction model would only correspond to the specific claim that it is more difficult *than you would already expect* to be in group ab knowing that someone is in group a and group b —or, alternatively, that is easier. The linear model claims that it is exactly the same as $a + b$.

What does the t -statistic on `female` now correspond to? It tests whether there is a gender effect among union members. Our baseline model only directly estimated this for non-union members.

Note that we could get both in one go with simply `reg wage4 ib1.female##ib1.union;` in that case, the model we are testing is $W = \beta_0 + \beta_1 N + \beta_2 M + \beta_3 NM + \epsilon$. Watch out for the signs, though; the direction of the effect depends on the baseline. Our non-intercept coefficients test, respectively, whether women get a union benefit, whether there is a gender effect among union members, and whether non-union men's mean is different to the additive effects of the first two (you might notice in this simple case that the interaction effect's t -statistic does not change across specifications).

1.1.2 Testing linear combinations

We can also get this information from many other approaches. Gordon mentions two: the F -test and tests of linear combinations. However, here, I will focus on using `margins` since that command is, in general, very useful. It is also shorter and simpler to use. Simply follow the syntax below (`dydx(x1)` refers to the partial derivative of our estimated regression function with respect to any `x1`, and we simply calculate it `at(x2=k)` for any value k). The underlying math is slightly different, but it comes to the same thing.

```
margins, dydx(union) at(female=1)
margins, dydx(female) at(union=1)
```

2 Categorical-continuous interactions

Another interesting model is one which involves an interaction between a dummy and a continuous variable: $Y = \beta_0 + \beta_1 D_1 + \beta_2 X_1 + \beta_3 D_1 X_1 + \epsilon$. What does this model say? Well, to work that out, let's start with something simpler.

You might now remember—or realize for the first time; that's OK!—that the simpler model $Y = \beta_0 + \beta_1 D_1 + \beta_2 X_1 + \epsilon$ gave us *only* lines with different intercepts: $Y =$

$$\begin{cases} \beta_0 + \beta_2 X_1; D_1 = 0 \\ \beta_0 + \beta_1 + \beta_2 X_1; D_1 = 1 \end{cases}$$

Let's look at an example with the CPS again. **Estimate a regression of wages on the female dummy and age without an interaction. Try to write the conditional regression equation for each group.**

```
reg wage4 i.female age
qui sum age
qui margins, at(female=(0 1) age = (`r(min)')(5)`r(max)'))
marginsplot, xdimension(age)
```

What does this model tell us? This simpler model tells us that wages are a function of someone's sex and their age, and it might seem almost dull; being female predicts lower earnings of \$4.75 an hour on average, one receives an estimated \$0.23 extra an hour for each year that they age, and a 0-year-old man would receive 18.60 an hour (hmmm...).

Note again that in this model, the effect of both of our variables is *constant*. They have *nothing* to do with any other variable. To see this, just compare “apples to apples”: make a prediction for a man with age $A = a$ and then “switch on” the female dummy; then, do this again but with age $A = a + 5$. The effect of sex will not change. Then, do the same but hold sex constant and examine the difference in predicted values between ages $A = a$ and $A = a + 5$, and then compare $A = a + 10$ and $A = a + 5$. In both cases, *neither* the value of the variable we hold constant *nor* the level of the variable we’re examining affects changes of equal size in the predictor. You can just see this in the `marginsplot` results.⁴

Perhaps what you are thinking is “that doesn’t sound very appealing as a model”. A more interesting model is one which involves an interaction between a dummy and a continuous variable: $Y = \beta_0 + \beta_1 D_1 + \beta_2 X_1 + \beta_3 D_1 X_1 + \epsilon$. We can plug in values that correspond to two groups $D_1 = 0, D_1 = 1$ since every person in the sample must be in one group.⁵ Then, we essentially have *two* models folded into one: $Y =$

$$\begin{cases} \beta_0 + \beta_2 X_1; D_1 = 0 \\ (\beta_0 + \beta_1) + A(\beta_2 + \beta_3); D_1 = 1 \end{cases}$$

You might notice that I did some creative grouping on the second equation—why? The point is to show that the two models might have essentially different parameters (the β s, the “control knobs” of the model) and simply happen to involve the same variables, D_1, X_1 (we’ll see later on that this model is the simplest case of a fully-interacted model). We likely have two intersecting lines: different slopes and intercepts. Let’s see that.

```
reg wage4 i.female##c.age
// note the interaction syntax here
qui sum age
margins, at(age = (`r(min)')(5)`r(max)')) over(nonwhite)
marginsplot
```

What does this model tell us? It’s maybe easier to write out the actual equations (with estimated, not known, parameters).

$$\begin{cases} 15.75 + 0.30A; F = 0 \\ (1.05 + 15.75) + A(0.30 - 0.14); F = 1 \end{cases}$$

Interestingly, we now also have *many different* sex effects; in the model without an interaction above, we had one. Now let’s again imagine predicting a value for someone in group $F = 0$ at some value $A = a$ and then “switch on” the dummy, i.e. change F to 1. Then we have $\mathbb{E}[Y|F = 0, A = a] = \beta_0 + \beta_2 a; \mathbb{E}[Y|F = 1, A = a] = (\beta_0 + \beta_1) + a(\beta_2 + \beta_3)$. The difference is $\beta_1 + a\beta_3$, i.e. the difference associated with “having been female” *is* now a function of some of the other variables in the model.

⁴If you want to be formal about it, you could also do this with expectations. Even simpler is to just take a partial derivative of the population regression function $Y : \frac{\partial Y}{\partial D} = \beta_1$. Since there is no D attached to any other variable, no other variable is involved when we just toggle D alone.

⁵Even if we have more dummies to represent a k -category polytomous variable, this always remains true of every individual dummy needed.

We can use `margins` to let us see that. Now we'll directly plot the effect of `female` at different values of `age`.

```
reg wage4 i.female##c.age
qui sum age
margins, dydx(female) at(age = (`r(min)')(5)`r(max)')
marginsplot
```

2.1 Centering and conditional effects

One interesting way to make our model more interpretable is to center the continuous predictor. This has the usual benefit that our intercept is now easier to interpret: it is the effect of being male (or, generally, in group $D_1 = 1$) *at the mean of our continuous predictor* (here, age). **Try that now with age.** Note that since we do not have simply *one* gender effect any more, our estimate of the intercept will change. However, we can actually use `margins` to test a specific effect. **Try centering age and then re-estimating the model. What do you notice?** `marginsplot` should help you interpret the results.

We can once again use `margins` to test conditional effects not shown in the model output. Notice that we are actually doing the same thing as before, although it is not obvious: when we were previously curious about what women's union effect was or what the gender effect was for union members, we were testing whether there was an effect of one variable, holding the other constant at some specific level. We can do the same now.

To take a simple example, let's reproduce our centered results using our original equation. We'll get a point estimate of the gender effect for individuals with the mean age.

```
reg wage4 i.female##c.age
qui sum age
margins, dydx(female) at(age=`r(mean)')
```

Perhaps we also want to test what the effect of gender on wages is at the start of working life and at the end (so, let's say, 18 and 67). **Give that a go.** If you like, you can show that this is the same as recentering the model at those points.

We can also go the other way and examine the effect of age across gender. Here, since we only have two values of gender, there is less to do, but we can still do it. **Verify that this gives the same result as re-estimating the model with women as the reference category using factor variable notation.**

```
reg wage4 i.female##c.age
qui sum age
margins, dydx(age) at(female=1)
```

3 The (so-called) Chow test

In cases where our model is *fully-interacted*, which is usually taken to mean specifically that we have a categorical variable that interacts with each of our other predictors, we can make clever use of the F -test to tell whether the effects of some relatively neutral “merit-type” variables (e.g. job tenure or age) vary by categorical variables on which people are often discriminated against—often corresponding to a natural and interesting sociological

hypothesis.⁶ Recall that we can always treat a partial F -test as just a ratio of the increment in SS_M between reduce and full to the full model's mean square error.

```
reg wage4 educ92 age
local SSm_r = e(mss)
reg wage4 i.female##c.educ92 i.female##c.age
local SSm_f = e(mss)
local SSm_incr = `SSm_f' - `SSm_r'
local MSE_f = e(rmse)^2
local df_diff = 3
local num = `SSm_incr' / `df_diff'
local F = `num' / `MSE_f'
di `F'
```

Now, if you don't like that approach, you can also go to the hassle of making product terms.

```
gen fem_ed = female*educ92
gen fem_age = female*age
reg wage4 female age educ92 fem_ed fem_age
test female fem_ed fem_age
```

However, there should probably be a faster way to do this, and there is.

```
contrast i.female#c.educ92 i.female#c.age, overall
```

4 Continuous-continuous interactions

Continuous-continuous interactions can be difficult to interpret. We'll discuss in class.

⁶If this interests you, I have a lot more to say about it; ask me for my notes about Oaxaca-Blinder decomposition. It is the most important method for studying discrimination and relies heavily on the categorical-continuous interaction approach to modeling and estimation.